

Semiparametric clustered overdispersed multinomial goodness-of-fit of log-linear models*

Alonso-Revenga, J. M.¹, Martín, N.^{2†}, Pardo, L.³

¹Department of Statistics and O.R. III, Complutense University of Madrid, Spain

²Department of Statistics and O.R. II, Complutense University of Madrid, Spain

³Department of Statistics and O.R. I, Complutense University of Madrid, Spain

September 26, 2016

Abstract

Traditionally, the Dirichlet-multinomial distribution has been recognized as a key model for contingency tables generated by cluster sampling schemes. There are, however, other possible distributions appropriate for these contingency tables. This paper introduces new test-statistics capable to test log-linear modeling hypotheses with no distributional specification, when the individuals of the clusters are possibly homogeneously correlated. The estimator for the intraclass correlation coefficient proposed in Alonso-Revenga et al. (2016), valid for different cluster sizes, plays a crucial role in the construction of the goodness-of-fit test-statistic.

Keywords: Clustered Multinomial Data; Consistent Intraclass Correlation Estimator; Log-linear model; Overdispersion; Quasi Minimum Divergence Estimator.

1 Introduction

In studies of frequency data, often the observations are organized in clusters. For clustered frequency data the classical statistical procedures are not longer valid. For example, in a study of hospitalized pairs of siblings, it is desired to study whether gender has any influence in schizophrenic diagnosis. Since the two outcomes of every pair of siblings (a cluster) are correlated for all the N pairs of siblings, the assumption of independence of all the $2N$ observations is violated and the classical independence test of two categorical variables, gender and schizophrenic diagnosis, is in principle useless. The same problem of invalidity of the classical chi-square and likelihood ratio tests are presented with any statistical model used for clustered frequencies.

Frequency data cross-classified according to K variables, (X_1, \dots, X_K) , having X_k categories $1, 2, \dots, I_k$, $k = 1, \dots, K$, are the so-called K -way contingency tables with $M = I_1 \times I_2 \times \dots \times I_K$ cells. In order to

*This paper was supported by the Spanish Grants MTM2015-67057 and ECO2015-66593 from Ministerio de Economía and Competitividad.

[†]Corresponding author, E-mail: nimartin@ucm.es.

clarify the concepts and notation we will focus our interest only on $K = 2$ variables, (X_1, X_2) , with I and J categories respectively, i.e. it has $M = I \times J$ cells denoted by pairs lexicographically ordered as

$$\Omega = \{(1, 1), (1, 2), \dots, (1, J), \dots, (I, 1), (I, 2), \dots, (I, J)\},$$

but it is possible to extend easily the same idea to K variables. The bidimensional random variable associated with the ℓ -th cluster of size n_ℓ , $\ell = 1, \dots, N$, being N the number of clusters, is denoted as

$$(X_{1,h}^{(\ell)}, X_{2,h}^{(\ell)}), \quad \ell = 1, \dots, N, \quad h = 1, \dots, n_\ell.$$

Let

$$I_S(X_1, X_2) = \begin{cases} 1, & \text{if } (X_1, X_2) \in S \\ 0, & \text{if } (X_1, X_2) \notin S \end{cases}$$

denote an indicator function of $S \subset \Omega$. Taking into account the total count associated with cell (i, j) is

$$Y_{ij}^{(\ell)} = \sum_{h=1}^{n_\ell} I_{\{(i,j)\}}(X_{1,h}^{(\ell)}, X_{2,h}^{(\ell)}), \quad \ell = 1, \dots, N, \quad (1.1)$$

the ℓ -th two-way frequency table in vector notation is given

$$\mathbf{Y}^{(\ell)} = (Y_{11}^{(\ell)}, \dots, Y_{1J}^{(\ell)}, \dots, Y_{I1}^{(\ell)}, \dots, Y_{IJ}^{(\ell)})^T, \quad \ell = 1, \dots, N,$$

where “ T ” denotes the transpose of a vector or matrix. In what follows, it is assumed an homogeneous probability for each individual felt in cell (i, j) of the ℓ -th cluster

$$p_{ij}(\boldsymbol{\theta}) = \Pr(X_{1,h}^{(\ell)} = i, X_{2,h}^{(\ell)} = j), \quad \ell = 1, \dots, N, \quad h = 1, \dots, n_\ell,$$

whose expression depends on an unknown M_0 -dimensional parameter vector

$$\boldsymbol{\theta} = (\theta_1, \dots, \theta_{M_0})^T \in \mathbb{R}^{M_0},$$

in terms of a log-linear model

$$\mathbf{p}(\boldsymbol{\theta}) = \frac{\exp\{\mathbf{W}\boldsymbol{\theta}\}}{\mathbf{1}_M^T \exp\{\mathbf{W}\boldsymbol{\theta}\}}, \quad (1.2)$$

where $M_0 < M - 1$,

$$\mathbf{p}(\boldsymbol{\theta}) = (p_{11}(\boldsymbol{\theta}), \dots, p_{1J}(\boldsymbol{\theta}), \dots, p_{I1}(\boldsymbol{\theta}), \dots, p_{IJ}(\boldsymbol{\theta}))^T \quad (1.3)$$

and the design matrix, \mathbf{W} , is a full rank matrix, with column vectors linearly independent with respect to the M -dimensional vector of 1's, $\mathbf{1}_M = (1, \dots, 1)^T$.

Under common correlation model for any pair of individuals h and s ($h, s = 1, \dots, n_\ell$, $h \neq s$) of any cluster $\ell = 1, \dots, N$, the intraclass correlation coefficient is defined as

$$\begin{aligned} \rho_{ij}^2 &= \text{Cor}[I_{\{(i,j)\}}(X_{1,h}^{(\ell)}, X_{2,h}^{(\ell)}), I_{\{(i,j)\}}(X_{1,s}^{(\ell)}, X_{2,s}^{(\ell)})] \\ &= \frac{E[I_{\{(i,j)\}}(X_{1,h}^{(\ell)}, X_{2,h}^{(\ell)}) I_{\{(i,j)\}}(X_{1,s}^{(\ell)}, X_{2,s}^{(\ell)})] - E[I_{\{(i,j)\}}(X_{1,h}^{(\ell)}, X_{2,h}^{(\ell)})] E[I_{\{(i,j)\}}(X_{1,s}^{(\ell)}, X_{2,s}^{(\ell)})]}{\sqrt{\text{Var}(I_{\{(i,j)\}}(X_{1,h}^{(\ell)}, X_{2,h}^{(\ell)})) \text{Var}(I_{\{(i,j)\}}(X_{1,s}^{(\ell)}, X_{2,s}^{(\ell)}))}} \\ &= \frac{\Pr(X_{1,h}^{(\ell)} = X_{1,s}^{(\ell)} = i, X_{2,h}^{(\ell)} = X_{2,s}^{(\ell)} = j) - p_{ij}^2(\boldsymbol{\theta})}{p_{ij}(\boldsymbol{\theta}) (1 - p_{ij}(\boldsymbol{\theta}))}, \quad \ell = 1, \dots, N, \quad h, s = 1, \dots, n_\ell, \quad h \neq s \end{aligned}$$

(see Eldridge et al. (2009), for more details). In correlated clustered overdispersed multinomial frequency data, in case of having homogeneous intracluster correlation cell by cell, $\rho^2 = \rho_{ij}^2$, $i = 1, \dots, I$, $j = 1, \dots, J$ and for this case, taking into account (1.1) and

$$\begin{aligned} E[I_{\{(i,j)\}}(X_{1,h}^{(\ell)}, X_{2,h}^{(\ell)})] &= p_{ij}(\boldsymbol{\theta}), \\ \text{Cov}[I_{\{(i,j)\}}(X_{1,h}^{(\ell)}, X_{2,h}^{(\ell)}), I_{\{(i,j)\}}(X_{1,s}^{(\ell)}, X_{2,s}^{(\ell)})] &= \begin{cases} \rho^2 p_{ij}(\boldsymbol{\theta}) (1 - p_{ij}(\boldsymbol{\theta})), & h \neq s \\ p_{ij}(\boldsymbol{\theta}) (1 - p_{ij}(\boldsymbol{\theta})), & h = s \end{cases}, \end{aligned}$$

it is proven that

$$E[\mathbf{Y}^{(\ell)}] = n_\ell \mathbf{p}(\boldsymbol{\theta}) \quad \text{and} \quad \text{Var}[\mathbf{Y}^{(\ell)}] = \vartheta_{n_\ell} n_\ell \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta})}, \quad (1.4)$$

where

$$\vartheta_{n_\ell} = 1 + (n_\ell - 1)\rho^2, \quad (1.5)$$

is referred to as “design effect” associated with the ℓ -th cluster,

$$\boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta})} = \mathbf{D}_{\mathbf{p}(\boldsymbol{\theta})} - \mathbf{p}(\boldsymbol{\theta})\mathbf{p}^T(\boldsymbol{\theta}), \quad (1.6)$$

and $\mathbf{D}_{\mathbf{p}(\boldsymbol{\theta})}$ is the diagonal matrix of $\mathbf{p}(\boldsymbol{\theta})$. Since $\text{Var}[Y_{ij}^{(\ell)}] > 0$, it holds $\vartheta_{n_\ell} = 1 + (n_\ell - 1)\rho^2 > 0$ for $\ell = 1, \dots, N$ and thus $\rho^2 > -1/(\max\{n_\ell\}_{\ell=1}^N - 1)$, but in practice it is assumed that $\rho^2 \geq 0$. This is just the reason why these models are termed “overdispersed models”. In particular, for $\rho^2 = 0$ all the frequency tables are multinomial.

Correlated clustered multinomial frequency data have been dealt in the statistical literature since many years ago through two different approaches. Following Choi and McHugh (1989), the design-based approach provides inferences with respect to the sampling distribution of estimates over repetitions of the same design. The works of Fellegi (1980), Holt et al. (1980), Rao and Scott (1981, 1984), Bedrick (1983), Landis et al (1984), Koch et al. (1975), Fay (1985), as well as references therein are good examples of this approach. On the other hand, Altham (1976), Cohen (1976), Brier (1980), Fienberg (1979), Menéndez et al. (1995, 1996) postulate a probability distribution to model the sample data. Dirichlet-multinomial is, historically, the first suitable distribution to modelize homogeneously correlated clustered overdispersed multinomial frequency with a fixed cluster size (see Mosimann, 1962). Later, Cohen (1976) and Altham (1976) proposed the n -inflated distribution and more recently, Morel and Nagaraj (1993) proposed the random-clumped distribution. The zero-inflated binomial distribution falls also inside this family of homogeneously correlated clustered overdispersed multinomial frequency data. Details about these distributions can be found in Alonso-Revenge et al. (2016). In the current paper and in Alonso-Revenge et al. (2016) a third approach is presented, different from the previous ones, based on the sole knowledge of the vector mean and the variance-covariance matrix of the distribution, given in (1.4), associated with the generator of the sample data. For log-linear modeling no distribution assumption is required if the quasi minimum ϕ -divergence estimators are used. In the following we shall assume that the data are generated by a population verifying (1.4). One of the strengths of this methodology, is that the proposed consistent estimator for ρ^2 is semi-parametric and it exhibits by far a better behavior with regard to the mean square error (MSE) in comparison with the existing estimation method, which is fully non-parametric. This kind of estimators are specially appealing for improving the behavior of the existing goodness-of-fit tests for log-linear models, with regard to the exact sizes and powers. The second strength of this methodology, is the flexibility in being applicable for different cluster sizes.

For the *semiparametric clustered overdispersed multinomial goodness-of-fit of log-linear models*, the interest lays on testing whether it holds a particular log-linear model

$$H_0 : \mathbf{p}(\boldsymbol{\theta}) = \frac{\exp\{\mathbf{W}\boldsymbol{\theta}\}}{\mathbf{1}_M^T \exp\{\mathbf{W}\boldsymbol{\theta}\}} \quad \text{vs.} \quad H_1 : \mathbf{p}(\boldsymbol{\theta}) \neq \frac{\exp\{\mathbf{W}\boldsymbol{\theta}\}}{\mathbf{1}_M^T \exp\{\mathbf{W}\boldsymbol{\theta}\}}. \quad (1.7)$$

2 Asymptotic Goodness-Of-Fit (GOF) test-statistics for equal cluster sizes

For the frequency tables and the probability vectors, a single index notation is preferred, since it covers any value, K , for the dimension of the contingency table. This means that the probability vector

$$\mathbf{p}(\boldsymbol{\theta}) = (p_1(\boldsymbol{\theta}), \dots, p_M(\boldsymbol{\theta}))^T,$$

and the ℓ -th frequency table

$$\mathbf{Y}^{(\ell)} = (Y_1^{(\ell)}, \dots, Y_M^{(\ell)})^T, \quad \ell = 1, \dots, N, \quad (2.1)$$

are valid to represent double index elements ordered as (1.3) when $K = 2$ ($M = IJ$), as well as to generalize for any value of K when the K -tuples are lexicographically ordered ($M = \prod_{k=1}^K I_k$). The M -dimensional vector obtained from collapsing the whole data, $\mathbf{Y}^{(\ell)}$, $\ell = 1, \dots, N$, is denoted by

$$\mathbf{Y} = \sum_{\ell=1}^N \mathbf{Y}^{(\ell)}$$

and MN -dimensional vector which gathers the whole data, $\mathbf{Y}^{(\ell)}$, $\ell = 1, \dots, N$, by

$$\tilde{\mathbf{Y}} = (\mathbf{Y}^{(1)T}, \dots, \mathbf{Y}^{(N)T})^T.$$

In this section, a family of GOF test-statistics for testing (1.7) with equal cluster sizes is introduced. In the following section the case of unequal cluster sizes is treated. Some preliminary results related to the estimators of the probability vector, derived in Alonso-Revilla et al. (2016), are first introduced. The non-parametric estimator of $\mathbf{p}(\boldsymbol{\theta})$, based on N clusters of sizes $n_\ell = n$, $\ell = 1, \dots, N$, is the M -dimensional vector of relative frequencies obtained collapsing the N frequency tables $\mathbf{Y}^{(\ell)}$, $\ell = 1, \dots, N$,

$$\hat{\mathbf{p}} = \frac{1}{nN} \mathbf{Y} = \frac{1}{nN} \sum_{\ell=1}^N \mathbf{Y}^{(\ell)} = \frac{1}{N} \sum_{\ell=1}^N \hat{\mathbf{p}}^{(\ell)},$$

where $\hat{\mathbf{p}}^{(\ell)} = \frac{1}{n} \mathbf{Y}^{(\ell)}$ represents the non-parametric estimator of $\mathbf{p}(\boldsymbol{\theta})$ based exclusively on the ℓ -th cluster.

Based on the collapsed table, \mathbf{Y} , the quasi minimum ϕ -divergence estimator (QM ϕ E) of $\boldsymbol{\theta}$ in (1.2) is defined as

$$\hat{\boldsymbol{\theta}}_\phi = \hat{\boldsymbol{\theta}}_\phi(\mathbf{Y}) = \arg \min_{\boldsymbol{\theta} \in \Theta} d_\phi(\hat{\mathbf{p}}, \mathbf{p}(\boldsymbol{\theta})),$$

where $\phi(x)$ is a convex function, $x > 0$, such that at $x = 1$, $\phi(1) = 0$, $\phi'(1) = 0$, $\phi''(1) > 0$, at $x = 0$, $0\phi(0/0) = 0$, $0\phi(p/0) = \lim_{u \rightarrow \infty} p\phi(u)/u$, and

$$d_\phi(\hat{\mathbf{p}}, \mathbf{p}(\boldsymbol{\theta})) = \sum_{r=1}^M p_r(\boldsymbol{\theta}) \phi\left(\frac{\hat{p}_r}{p_r(\boldsymbol{\theta})}\right) \quad (2.2)$$

is the ϕ -divergence between the probability vectors $\hat{\mathbf{p}}$ and $\mathbf{p}(\boldsymbol{\theta})$. For more details about ϕ -divergence measures see Cressie and Pardo (2002) and Pardo (2006).

The *quasi-maximum likelihood estimator* (QMLE) of $\boldsymbol{\theta}$, denoted by $\hat{\boldsymbol{\theta}}$, is a particular case of the QM ϕ E by replacing the ϕ -divergence by the Kullback divergence between the probability vectors $\hat{\mathbf{p}}$ and $\mathbf{p}(\boldsymbol{\theta})$, i.e.,

$$\begin{aligned} \hat{\boldsymbol{\theta}} &= \hat{\boldsymbol{\theta}}(\mathbf{Y}) = \arg \min_{\boldsymbol{\theta} \in \Theta} d_{Kullback}(\hat{\mathbf{p}}, \mathbf{p}(\boldsymbol{\theta})), \\ d_{Kullback}(\hat{\mathbf{p}}, \mathbf{p}(\boldsymbol{\theta})) &= \sum_{r=1}^M \hat{p}_r \log \frac{\hat{p}_r}{p_r(\boldsymbol{\theta})}, \end{aligned}$$

or equivalently $\hat{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}(\mathbf{Y}) = \arg \min_{\boldsymbol{\theta} \in \Theta} d_\phi(\hat{\mathbf{p}}, \mathbf{p}(\boldsymbol{\theta}))$, with $\phi(x) = x \log x - x + 1$. Since the QM ϕ Es are invariant estimators,

$$\mathbf{p}(\hat{\boldsymbol{\theta}}_\phi) = \frac{\exp\{\mathbf{W}\hat{\boldsymbol{\theta}}_\phi\}}{\mathbf{1}_M^T \exp\{\mathbf{W}\hat{\boldsymbol{\theta}}_\phi\}}$$

is the QM ϕ Es of $\mathbf{p}(\boldsymbol{\theta})$.

Theorem 2.1 *The asymptotic distribution of the difference between the non-parametric estimator and the QM ϕ E of $\mathbf{p}(\boldsymbol{\theta})$, with N clusters of size n , is*

$$\sqrt{N}(\hat{\mathbf{p}} - \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_2})) \xrightarrow[N \rightarrow \infty]{\mathcal{L}} \mathcal{N}(\mathbf{0}_M, \frac{\vartheta_n}{n}(\boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} - \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} \mathbf{W} (\mathbf{W}^T \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} \mathbf{W})^{-1} \mathbf{W}^T \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)})),$$

where $\boldsymbol{\theta}_0$ is the unknown true value of $\boldsymbol{\theta}$.

Proof. By following (A.5) and (A.3) in the proof of Theorem 2.2 of Alonso-Revilla et al. (2016, Section A.3), it holds

$$\sqrt{N}(\mathbf{p}(\hat{\boldsymbol{\theta}}_\phi) - \mathbf{p}(\boldsymbol{\theta}_0)) = \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} \mathbf{W} \sqrt{N}(\hat{\boldsymbol{\theta}}_\phi - \boldsymbol{\theta}_0) + o_p(\mathbf{1}_M),$$

and

$$\sqrt{N}(\hat{\boldsymbol{\theta}}_\phi - \boldsymbol{\theta}_0) = \mathbf{D}_{\mathbf{p}(\boldsymbol{\theta}_0)}^{-1/2} (\mathbf{A}^T(\boldsymbol{\theta}_0) \mathbf{A}(\boldsymbol{\theta}_0))^{-1} \mathbf{A}^T(\boldsymbol{\theta}_0) \sqrt{N}(\hat{\mathbf{p}} - \mathbf{p}(\boldsymbol{\theta}_0)) + o_p(\mathbf{1}_{M_0}),$$

where

$$\mathbf{A}(\boldsymbol{\theta}_0) = \mathbf{D}_{\mathbf{p}(\boldsymbol{\theta}_0)}^{-1/2} \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} \mathbf{W}.$$

Plugging $\sqrt{N}(\hat{\boldsymbol{\theta}}_\phi - \boldsymbol{\theta}_0)$ into the expression of $\sqrt{N}(\mathbf{p}(\hat{\boldsymbol{\theta}}_\phi) - \mathbf{p}(\boldsymbol{\theta}_0))$ we get

$$\sqrt{N}(\mathbf{p}(\hat{\boldsymbol{\theta}}_\phi) - \mathbf{p}(\boldsymbol{\theta}_0)) = \mathbf{A}(\boldsymbol{\theta}_0) (\mathbf{A}^T(\boldsymbol{\theta}_0) \mathbf{A}(\boldsymbol{\theta}_0))^{-1} \mathbf{A}^T(\boldsymbol{\theta}_0) \sqrt{N}(\hat{\mathbf{p}} - \mathbf{p}(\boldsymbol{\theta}_0)) + o_p(\mathbf{1}_M),$$

and subtracting the expressions on both sides of the equality to $\hat{\mathbf{p}} - \mathbf{p}(\boldsymbol{\theta}_0)$,

$$\sqrt{N}(\hat{\mathbf{p}} - \mathbf{p}(\hat{\boldsymbol{\theta}}_\phi)) = \left(\mathbf{I}_M - \mathbf{A}(\boldsymbol{\theta}_0) (\mathbf{A}^T(\boldsymbol{\theta}_0) \mathbf{A}(\boldsymbol{\theta}_0))^{-1} \mathbf{A}^T(\boldsymbol{\theta}_0) \right) \sqrt{N}(\hat{\mathbf{p}} - \mathbf{p}(\boldsymbol{\theta}_0)) + o_p(\mathbf{1}_M).$$

On the other hand, by applying the Central Limit Theorem

$$\sqrt{N}(\hat{\mathbf{p}} - \mathbf{p}(\boldsymbol{\theta}_0)) \xrightarrow[N \rightarrow \infty]{\mathcal{L}} \mathcal{N}(\mathbf{0}_M, \frac{\vartheta_n}{n} \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)}) \quad (2.3)$$

(see Alonso-Revenge et al. (2016), eq. (3.1)), from which the asymptotic distribution of $\sqrt{N}(\hat{\mathbf{p}} - \mathbf{p}(\hat{\boldsymbol{\theta}}_\phi))$ is an M -dimensional central normal with variance-covariance matrix equal to

$$\begin{aligned} & \frac{\vartheta_n}{n} \left(\mathbf{I}_M - \mathbf{A}(\boldsymbol{\theta}_0) (\mathbf{A}^T(\boldsymbol{\theta}_0) \mathbf{A}(\boldsymbol{\theta}_0))^{-1} \mathbf{A}^T(\boldsymbol{\theta}_0) \right) \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} \left(\mathbf{I}_M - \mathbf{A}(\boldsymbol{\theta}_0) (\mathbf{A}^T(\boldsymbol{\theta}_0) \mathbf{A}(\boldsymbol{\theta}_0))^{-1} \mathbf{A}^T(\boldsymbol{\theta}_0) \right) \\ &= \frac{\vartheta_n}{n} (\boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} - \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} \mathbf{W} (\mathbf{W}^T \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} \mathbf{W})^{-1} \mathbf{W}^T \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)}). \end{aligned}$$

The last equality comes from $\boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} \mathbf{D}_{\mathbf{p}(\boldsymbol{\theta}_0)}^{-1} \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} = \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)}$ and $\mathbf{D}_{\mathbf{p}(\boldsymbol{\theta}_0)}^{-1} \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} \mathbf{A}(\boldsymbol{\theta}_0) = \mathbf{A}(\boldsymbol{\theta}_0)$. ■

The *semi-parametric estimator* of ϑ_n , via QMφEs, is

$$\tilde{\vartheta}_{n,N,\phi} = \frac{X^2(\tilde{\mathbf{Y}}, \hat{\boldsymbol{\theta}}_\phi)}{(N-1)(M-1)},$$

where

$$X^2(\tilde{\mathbf{Y}}, \hat{\boldsymbol{\theta}}_\phi) = \sum_{\ell=1}^N \left(\mathbf{Y}^{(\ell)} - n\hat{\mathbf{p}} \right)^T \frac{1}{n} \mathbf{D}_{\hat{\mathbf{p}}(\hat{\boldsymbol{\theta}}_\phi)}^{-1} \left(\mathbf{Y}^{(\ell)} - n\hat{\mathbf{p}} \right) = n \sum_{\ell=1}^N \sum_{r=1}^M \frac{(\hat{p}_r^{(\ell)} - \hat{p}_r)^2}{p_r(\hat{\boldsymbol{\theta}}_\phi)}. \quad (2.4)$$

Similarly, the semi-parametric estimator of ρ^2 , via QMφEs, is

$$\tilde{\rho}_{n,N,\phi}^2 = \frac{\tilde{\vartheta}_{n,N,\phi} - 1}{n - 1}.$$

Both, $\tilde{\vartheta}_{n,N,\phi}$ and $\tilde{\rho}_{n,N,\phi}^2$, are consistent estimators of ϑ and ρ^2 respectively.

Corollary 2.2 *The semiparametric clustered overdispersed chi-square GOF test-statistic, with N clusters of size n , has the following asymptotic distribution*

$$\frac{X^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_\phi)}{\tilde{\vartheta}_{n,N,\phi}} \xrightarrow[N \rightarrow \infty]{\mathcal{L}} \chi_{M-M_0-1}^2,$$

where

$$X^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_\phi) = nN(\hat{\mathbf{p}} - \mathbf{p}(\hat{\boldsymbol{\theta}}_\phi))^T \mathbf{D}_{\mathbf{p}(\hat{\boldsymbol{\theta}}_\phi)}^{-1} (\hat{\mathbf{p}} - \mathbf{p}(\hat{\boldsymbol{\theta}}_\phi)) = nN \sum_{r=1}^M \frac{(\hat{p}_r - p_r(\hat{\boldsymbol{\theta}}_\phi))^2}{p_r(\hat{\boldsymbol{\theta}}_\phi)}. \quad (2.5)$$

Proof.

$$\frac{X^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_\phi)}{\tilde{\vartheta}_{n,N,\phi}} = \mathbf{Q}^T \mathbf{Q}, \quad (2.6)$$

where

$$Q = \sqrt{\frac{n}{\tilde{\vartheta}_{n,N,\phi}}} D_{\mathbf{p}(\hat{\boldsymbol{\theta}}_\phi)}^{-1/2} \sqrt{N}(\hat{\mathbf{p}} - \mathbf{p}(\hat{\boldsymbol{\theta}}_\phi))$$

is an M -dimensional central normal with variance-covariance matrix equal to

$$\begin{aligned} \mathbf{V}(\boldsymbol{\theta}_0) &= D_{\mathbf{p}(\boldsymbol{\theta}_0)}^{-1} \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} - D_{\mathbf{p}(\hat{\boldsymbol{\theta}}_\phi)}^{-1/2} \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} \mathbf{W} (\mathbf{W}^T \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} \mathbf{W})^{-1} \mathbf{W}^T \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} D_{\mathbf{p}(\boldsymbol{\theta}_0)}^{-1/2} \\ &= D_{\mathbf{p}(\boldsymbol{\theta}_0)}^{-1} \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} - \mathbf{A}(\boldsymbol{\theta}_0) (\mathbf{A}^T(\boldsymbol{\theta}_0) \mathbf{A}(\boldsymbol{\theta}_0))^{-1} \mathbf{A}^T(\boldsymbol{\theta}_0), \end{aligned}$$

by applying the Slutsky's Theorem. The asymptotic distribution of a quadratic form, such as (2.6), with $\mathbf{V}(\boldsymbol{\theta}_0)$ being idempotent, is a chi-square distribution with degrees of freedom equal to the rank of $\mathbf{V}(\boldsymbol{\theta}_0)$. The idempotence of $\mathbf{V}(\boldsymbol{\theta}_0)$ is proven with similar arguments given to obtain the variance-covariance matrix at the end of the proof of Theorem 2.1. Finally, taking into account that for idempotent matrices rank and trace are equivalent, and by properties of the trace of the product of two matrices, it holds

$$\begin{aligned} \text{rank}(\mathbf{V}(\boldsymbol{\theta}_0)) &= \text{trace}(\mathbf{V}(\boldsymbol{\theta}_0)) \\ &= \text{trace}(D_{\mathbf{p}(\boldsymbol{\theta}_0)}^{-1} \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)}) - \text{trace}(\mathbf{A}(\boldsymbol{\theta}_0) (\mathbf{A}^T(\boldsymbol{\theta}_0) \mathbf{A}(\boldsymbol{\theta}_0))^{-1} \mathbf{A}^T(\boldsymbol{\theta}_0)) \\ &= \text{trace}(\mathbf{I}_M - \mathbf{p}(\boldsymbol{\theta}_0)) - \text{trace}(\mathbf{I}_{M_0}) \\ &= (M - 1) - M_0. \end{aligned}$$

■

Remark 2.3 *The chi-square statistics given in (2.4) and (2.5) need some clarifications, since under the same terminology arise totally different ideas. While $X^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_\phi)$ is part of a GOF test-statistic, $X^2(\tilde{\mathbf{Y}}, \hat{\boldsymbol{\theta}}_\phi)$ is part of an estimator constructed through the trace of the quasi-variance-covariance matrix of $\tilde{\mathbf{Y}}$ (see more details in Alonso-Revenge et al. (2016)). Structurally, $X^2(\tilde{\mathbf{Y}}, \hat{\boldsymbol{\theta}}_\phi)$ is quite different from the usual chi-square test-statistics, since the total number of cells, NM , depends on N , which increases to infinity.*

The ϕ -divergence measures permit to construct either estimators as well as test-statistics. Both of them do not need to be the same, for example in the usual chi-square test-statistic with QMLEs, $X^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}) = 2nN d_{\phi_1}(\hat{\mathbf{p}}, \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_2}))$, where $\phi_1(x) = \frac{1}{2}(x^2 - 1)$ and $\phi_2(x) = x \log x - x + 1$. In what is to follow, notation ϕ_1 and ϕ_2 are used to distinguish the ϕ function of the ϕ -divergences.

Theorem 2.4 *The semiparametric clustered overdispersed divergence based GOF test-statistic, with N clusters of size n , has the following asymptotic distribution*

$$\frac{T^{\phi_1}(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_2})}{\tilde{\vartheta}_{n,N,\phi_2}} \xrightarrow[N \rightarrow \infty]{\mathcal{L}} \chi_{M-M_0-1}^2,$$

where

$$T^{\phi_1}(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_2}) = \frac{2nN}{\phi_1''(1)} d_{\phi_1}(\hat{\mathbf{p}}, \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_2})) \quad (2.7)$$

and

$$d_{\phi_1}(\hat{\mathbf{p}}, \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_2})) = \sum_{r=1}^M p_r(\hat{\boldsymbol{\theta}}_{\phi_2}) \phi_1\left(\frac{\hat{p}_r}{p_r(\hat{\boldsymbol{\theta}}_{\phi_2})}\right).$$

Proof. A second order Taylor expansion of $d_{\phi_1}(\hat{\mathbf{p}}, \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_2}))$ around $(\mathbf{p}(\boldsymbol{\theta}_0), \mathbf{p}(\boldsymbol{\theta}_0))$ needs derivatives of first and second order. Since

$$d_{\phi_1}(\mathbf{p}, \mathbf{q}) = \sum_{j=1}^M q_j \phi_1\left(\frac{p_j}{q_j}\right),$$

and $\phi_1(1) = \phi_1'(1) = 0$, the first order derivatives of the Taylor expansion are cancelled. The second order derivatives yields

$$\begin{aligned} \left. \frac{\partial^2 d_{\phi_1}(\mathbf{p}, \mathbf{p}(\boldsymbol{\theta}_0))}{\partial \mathbf{p} \partial \mathbf{p}^T} \right|_{\mathbf{p}=\mathbf{p}(\boldsymbol{\theta}_0)} &= \left. \frac{\partial^2 d_{\phi_1}(\mathbf{p}(\boldsymbol{\theta}_0), \mathbf{q})}{\partial \mathbf{q} \partial \mathbf{q}^T} \right|_{\mathbf{q}=\mathbf{p}(\boldsymbol{\theta}_0)} = - \left. \frac{\partial^2 d_{\phi_1}(\mathbf{p}, \mathbf{q})}{\partial \mathbf{q} \partial \mathbf{p}^T} \right|_{\mathbf{p}=\mathbf{p}(\boldsymbol{\theta}_0), \mathbf{q}=\mathbf{p}(\boldsymbol{\theta}_0)} \\ &= \phi_1''(1) \mathbf{D}_{\mathbf{p}(\boldsymbol{\theta}_0)}^{-1}. \end{aligned}$$

Hence,

$$d_{\phi_1}(\hat{\mathbf{p}}, \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_2})) = \frac{1}{2} \phi_1''(1) (\hat{\mathbf{p}} - \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_2}))^T \mathbf{D}_{\mathbf{p}(\boldsymbol{\theta}_0)}^{-1} (\hat{\mathbf{p}} - \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_2})) + o(\|\hat{\mathbf{p}} - \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_2})\|^2).$$

Finally, since Theorem 2.1 $o(N\|\hat{\mathbf{p}} - \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_2})\|^2) = o_p(1)$, and thus

$$\frac{T^{\phi_1}(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_2})}{\tilde{\vartheta}_{n,N,\phi_2}} = \frac{1}{\tilde{\vartheta}_{n,N,\phi_2}} \frac{2Nn}{\phi_1''(1) \vartheta_n} d_{\phi_1}(\hat{\mathbf{p}}, \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_2})) = \frac{1}{\tilde{\vartheta}_{n,N,\phi_2}} \mathbf{Q}^T \mathbf{Q} + o_p(1) = \frac{X^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_2})}{\tilde{\vartheta}_{n,N,\phi_2}} + o_p(1),$$

which means that $T^{\phi_1}(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_2})/\tilde{\vartheta}_{n,N,\phi_2}$ and $X^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_2})/\tilde{\vartheta}_{n,N,\phi_2}$ have the same asymptotic distribution, $\chi_{M-M_0-1}^2$, according to Corollary 2.2. ■

The following result is a particular case of Theorem 2.4, with $\phi_1(x) = x \log x - x + 1$.

Corollary 2.5 *The semiparametric clustered overdispersed likelihood-ratio GOF test-statistic, with N clusters of size n , has the following asymptotic distribution*

$$\frac{G^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_2})}{\tilde{\vartheta}_{n,N,\phi_2}} \xrightarrow[N \rightarrow \infty]{\mathcal{L}} \chi_{M-M_0-1}^2,$$

where

$$G^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_2}) = nN \sum_{r=1}^M \hat{p}_r \log \frac{\hat{p}_r}{p_r(\hat{\boldsymbol{\theta}}_{\phi_2})}. \quad (2.8)$$

3 Asymptotic Goodness-Of-Fit (GOF) test-statistics for unequal cluster sizes

As introduction of this section a brief summary of the estimators given in Alonso-Revilla et al. (2016) is presented. Organizing the frequency tables associated with the clusters, according to their sizes, the double index in

$$\mathbf{Y}^{(g,\ell)} = (Y_1^{(g,\ell)}, \dots, Y_M^{(g,\ell)})^T, \quad \ell = 1, \dots, N,$$

denotes the g -th frequency table of size n_ℓ , $g = 1, \dots, G$. In this setting, the non-parametric estimator of $\mathbf{p}(\boldsymbol{\theta})$ is, according to Alonso-Revenge et al. (2016), equal to

$$\hat{\mathbf{p}} = \frac{1}{\sum_{h=1}^G n_h N_h} \mathbf{Y} = \frac{1}{\sum_{h=1}^G n_h N_h} \sum_{g=1}^G \sum_{\ell=1}^N \mathbf{Y}^{(g,\ell)} = \sum_{g=1}^G w_g \hat{\mathbf{p}}^{(g)},$$

where

$$\begin{aligned} \mathbf{Y} &= \sum_{g=1}^G \sum_{\ell=1}^N \mathbf{Y}^{(g,\ell)}, \\ w_g &= \frac{n_g N_g}{\sum_{h=1}^G n_h N_h}, \end{aligned} \tag{3.1}$$

and

$$\begin{aligned} \hat{\mathbf{p}}^{(g)} &= \frac{1}{n_g N} \sum_{\ell=1}^{N_g} \mathbf{Y}^{(g,\ell)} = \frac{1}{N} \sum_{\ell=1}^{N_g} \hat{\mathbf{p}}^{(g,\ell)}, \\ \hat{\mathbf{p}}^{(g,\ell)} &= \frac{1}{n_g} \mathbf{Y}^{(g,\ell)}. \end{aligned}$$

Let

$$\vartheta_{n^*} = 1 + \rho^2 (n^* - 1) \in (1, n^*], \tag{3.2}$$

be design effect with a sample size equal to

$$\begin{aligned} n^* &= \sum_{g=1}^G w_g^* n_g, \\ w_g^* &= \frac{N_g^* n_g}{\sum_{h=1}^G N_h^* n_h} > 0, \quad g = 1, \dots, G, \end{aligned}$$

and $N_g^* \in (0, 1]$ an unknown value such that

$$\frac{N_g}{N} \xrightarrow[N \rightarrow \infty]{P} N_g^*, \quad g = 1, \dots, G.$$

The *semi-parametric estimator* of ϑ_{n^*} , via QM ϕ Es, is

$$\tilde{\vartheta}_{\hat{n}^*, N, \phi} = \sum_{g=1}^G w_g \tilde{\vartheta}_{n_g, N_g, \phi}, \tag{3.3}$$

where w_g is (3.1),

$$\hat{n}^* = \sum_{g=1}^G w_g n_g,$$

$$N = \sum_{g=1}^G N_g,$$

$$\tilde{\vartheta}_{n_g, N_g, \phi} = \frac{X^2(\tilde{\mathbf{Y}}_g, \hat{\boldsymbol{\theta}}_\phi)}{(N_g - 1)(M - 1)},$$

$$X^2(\tilde{\mathbf{Y}}_g, \hat{\boldsymbol{\theta}}_\phi) = n_g \sum_{\ell=1}^{N_g} \sum_{r=1}^M \frac{(\hat{p}_r^{(\ell, g)} - \hat{p}_r^{(g)})^2}{p_r(\hat{\boldsymbol{\theta}}_\phi)} = n_g \sum_{r=1}^M \frac{1}{p_r(\hat{\boldsymbol{\theta}}_\phi)} \sum_{\ell=1}^{N_g} (\hat{p}_r^{(\ell, g)} - \hat{p}_r^{(g)})^2.$$

Similarly, the semi-parametric estimator of ρ^2 , via QM ϕ Es, is

$$\tilde{\rho}_{\hat{n}^*, N, \phi}^2 = \frac{\tilde{\vartheta}_{\hat{n}^*, N, \phi} - 1}{\hat{n}^* - 1}.$$

Both, $\tilde{\vartheta}_{\hat{n}^*, N, \phi}$ and $\tilde{\rho}_{\hat{n}^*, N, \phi}^2$, are consistent estimators of ϑ and ρ^2 respectively.

The following results are not explicitly proven since the same steps of the proof given in Section 2 are needed. However, a basic and different result is required in the place of (2.3), which is

$$\sqrt{N}(\hat{\mathbf{p}} - \mathbf{p}(\boldsymbol{\theta}_0)) \xrightarrow[N \rightarrow \infty]{\mathcal{L}} \mathcal{N}(\mathbf{0}_M, \frac{\vartheta_{n^*}}{\bar{n}} \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)}), \quad (3.4)$$

proven in Alonso-Revilla et al. (2016).

Theorem 3.1 *The asymptotic distribution of the difference between the non-parametric estimator and QM ϕ E of $\mathbf{p}(\boldsymbol{\theta})$, with G groups of clusters of size n_g , $g = 1, \dots, G$, is*

$$\sqrt{N}(\hat{\mathbf{p}} - \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_2})) \xrightarrow[N \rightarrow \infty]{\mathcal{L}} \mathcal{N}(\mathbf{0}_M, \frac{\vartheta_{n^*}}{\bar{n}} (\boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} - \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} \mathbf{W} (\mathbf{W}^T \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)} \mathbf{W})^{-1} \mathbf{W}^T \boldsymbol{\Sigma}_{\mathbf{p}(\boldsymbol{\theta}_0)})),$$

where $\boldsymbol{\theta}_0$ is the unknown true value of $\boldsymbol{\theta}$, ϑ_{n^*} is (3.2) and

$$\bar{n} = \sum_{g=1}^G N_g^* n_g.$$

Corollary 3.2 *The semiparametric clustered overdispersed chi-square GOF test-statistic, with G groups of clusters of size n_g , $g = 1, \dots, G$, has the following asymptotic distribution*

$$\frac{X^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_2})}{\tilde{\vartheta}_{\hat{n}^*, N, \phi_2}} \xrightarrow[N \rightarrow \infty]{\mathcal{L}} \chi_{M-M_0-1}^2,$$

where $\tilde{\vartheta}_{\hat{n}^*, N, \phi}$ is (3.3) and

$$X^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_2}) = \hat{n}N(\hat{\mathbf{p}} - \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_2}))^T \mathbf{D}_{\mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_2})}^{-1} (\hat{\mathbf{p}} - \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_2})) = \hat{n}N \sum_{r=1}^M \frac{(\hat{p}_r - p_r(\hat{\boldsymbol{\theta}}_{\phi_2}))^2}{p_r(\hat{\boldsymbol{\theta}}_{\phi_2})},$$

with

$$\hat{n}N = \sum_{g=1}^G N_g n_g. \quad (3.5)$$

Theorem 3.3 *The semiparametric clustered overdispersed divergence based GOF test-statistic, with G groups of clusters of size n_g , $g = 1, \dots, G$, has the following asymptotic distribution*

$$\frac{T^{\phi_1}(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_2})}{\tilde{\vartheta}_{\hat{n}^*, N, \phi_2}} \xrightarrow[N \rightarrow \infty]{\mathcal{L}} \chi_{M-M_0-1}^2,$$

where

$$T^{\phi_1}(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_2}) = \frac{2\hat{n}N}{\phi_1''(1)} d_{\phi_1}(\hat{\mathbf{p}}, \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_2})),$$

$\tilde{\vartheta}_{\hat{n}^*, N, \phi_2}$ is (3.3) and $\hat{n}N$ (3.5).

Corollary 3.4 *The semiparametric clustered overdispersed likelihood-ratio GOF test-statistic, with G groups of clusters of size n_g , $g = 1, \dots, G$, has the following asymptotic distribution*

$$\frac{G^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_2})}{\tilde{\vartheta}_{\hat{n}^*, N, \phi_2}} \xrightarrow[N \rightarrow \infty]{\mathcal{L}} \chi_{M-M_0-1}^2,$$

where

$$G^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_2}) = \hat{n}N \sum_{r=1}^M \hat{p}_r \log \frac{\hat{p}_r}{p_r(\hat{\boldsymbol{\theta}}_{\phi_2})},$$

$\tilde{\vartheta}_{\hat{n}^*, N, \phi_2}$ is (3.3) and $\hat{n}N$ (3.5).

Brier (1980) proposed using $G^2(\mathbf{Y}, \hat{\boldsymbol{\theta}})/\hat{\vartheta}_{\hat{n}^*, N}$ and $X^2(\mathbf{Y}, \hat{\boldsymbol{\theta}})/\hat{\vartheta}_{\hat{n}^*, N}$, with

$$\hat{\vartheta}_{\hat{n}^*, N} = \sum_{g=1}^G w_g \hat{\vartheta}_{n_g, N_g}, \quad (3.6)$$

where

$$\hat{\vartheta}_{n_g, N_g} = \frac{n_g}{(N_g - 1)(M - 1)} \sum_{r=1}^M \frac{1}{\hat{p}_r^{(g)}} \sum_{\ell=1}^{N_g} (\hat{p}_r^{(\ell, g)} - \hat{p}_r^{(g)})^2,$$

to be applied for the Dirichlet-multinomial distribution. In a similar way as Theorem 3.3, it is proven that

$$\frac{T^{\phi_1}(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_2})}{\hat{\vartheta}_{n_g, N_g}} \xrightarrow[N \rightarrow \infty]{\mathcal{L}} \chi_{M-M_0-1}^2. \quad (3.7)$$

4 Numerical example

From all the households located in $N = 20$ neighborhoods around Montevideo (Minnesota, US), some households were randomly selected: from $N_1 = 18$ neighborhoods $n_1 = 5$ houses were selected and from $N_2 = 2$ neighborhoods $n_2 = 3$ houses. The neighborhoods are grouped into class $g = 1$ or $g = 2$ depending on the selected number of houses (neighborhood or cluster size), $n_1 = 5$ and $n_2 = 3$ respectively. For the ℓ -th neighborhood ($\ell = 1, \dots, N_g$) of the g -th cluster size, in the s -th selected

home ($s = 1, \dots, n_g$), the family was questioned on two study interests: satisfaction with the housing in the neighborhood as a whole ($X_{1s}^{(g,\ell)}$), and satisfaction with their own home ($X_{2s}^{(g,\ell)}$). For both questions the responses were classified as unsatisfied (US), satisfied (S) or very satisfied (VS). In the sequel, we shall identify the aforementioned categories of the ordinal variables, $X_{11}^{(g,\ell)}$ and $X_{12}^{(g,\ell)}$, with numbers 1, 2, and 3: for example, (US, S) is associated with $(X_{11}^{(g,\ell)}, X_{12}^{(g,\ell)}) = (1, 2)$.

Under the null hypothesis of (1.7), a family's classification according to level of personal satisfaction is independent from its classification by level of community satisfaction. The corresponding log-linear model, $\log p_{ij}(\boldsymbol{\theta}) = u + \theta_{1(i)} + \theta_{2(j)}$, for $i = 1, \dots, I = 3$, $j = 1, \dots, J = 3$, has as design matrix and the unknown parameter vector

$$\mathbf{W} = \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 & -1 & -1 & -1 \\ 0 & 0 & 0 & 1 & 1 & 1 & -1 & -1 & -1 \\ 1 & 0 & -1 & 1 & 0 & -1 & 1 & 0 & -1 \\ 0 & 1 & -1 & 0 & 1 & -1 & 0 & 1 & -1 \end{pmatrix}^T \quad \text{and} \quad \boldsymbol{\theta} = (\theta_{1(1)}, \theta_{1(2)}, \theta_{2(1)}, \theta_{2(2)})^T.$$

The corresponding data, given in Table 4.1, are disaggregated based on the number of houses (g) and neighborhood identifications (ℓ) in 20 rows, having each $M = 9$ cells in lexicographical order. The $G = 2$ groups of clusters have respectively $n_1 = 5$ and $n_2 = 3$ families.

g	ℓ	$Y_{11}^{(g,\ell)}$	$Y_{12}^{(g,\ell)}$	$Y_{13}^{(g,\ell)}$	$Y_{21}^{(g,\ell)}$	$Y_{22}^{(g,\ell)}$	$Y_{23}^{(g,\ell)}$	$Y_{31}^{(g,\ell)}$	$Y_{32}^{(g,\ell)}$	$Y_{33}^{(g,\ell)}$
1	1	1	0	0	2	2	0	0	0	0
1	2	1	0	0	2	2	0	0	0	0
1	3	0	2	0	0	2	0	0	1	0
1	4	0	1	0	2	1	0	1	0	0
1	5	0	0	0	0	4	0	0	1	0
1	6	1	0	0	3	1	0	0	0	0
1	7	3	0	0	0	1	0	0	1	0
1	8	1	0	0	1	3	0	0	0	0
1	9	3	0	0	0	0	0	1	0	1
1	10	0	1	0	0	3	1	0	0	0
1	11	1	1	0	0	2	0	1	0	0
1	12	0	1	0	4	0	0	0	0	0
1	13	0	0	0	4	1	0	0	0	0
1	14	0	0	0	1	2	0	0	0	2
1	15	2	0	0	2	1	0	0	0	0
1	16	0	0	0	1	1	1	0	2	0
1	17	2	0	0	2	1	0	0	0	0
1	18	2	0	0	2	0	0	1	0	0
2	1	1	0	0	1	1	0	0	0	0
2	2	0	0	0	1	0	1	0	0	1

Table 4.1: Housing satisfaction in neighbourhoods of Montevideo (Brier, 1980).

For estimation and testing, the power divergence measures are considered, by restricting ϕ from the family of convex functions to the subfamily

$$\phi_\lambda(x) = \begin{cases} \frac{1}{\lambda(1+\lambda)} [x^{\lambda+1} - x - \lambda(x-1)], & \lambda \notin \{-1, 0\} \\ \lim_{v \rightarrow \lambda} \frac{1}{v(1+v)} [x^{v+1} - x - v(x-1)], & \lambda \in \{-1, 0\} \end{cases},$$

where $\lambda \in \mathbb{R}$ is a tuning parameter. The expression of (2.2) becomes

$$d_{\phi_\lambda}(\hat{\mathbf{p}}, \mathbf{p}(\boldsymbol{\theta})) = \begin{cases} \frac{1}{\lambda(\lambda+1)} \left(\sum_{r=1}^M \frac{\hat{p}_r^{\lambda+1}}{p_r^\lambda(\boldsymbol{\theta})} - 1 \right), & \lambda \notin \{-1, 0\} \\ d_{Kullback}(\mathbf{p}(\boldsymbol{\theta}), \hat{\mathbf{p}}), & \lambda = -1 \\ d_{Kullback}(\hat{\mathbf{p}}, \mathbf{p}(\boldsymbol{\theta})), & \lambda = 0 \end{cases},$$

in such a way that for each $\lambda \in \mathbb{R}$ a different divergence measure is obtained. The quasi minimum power-divergence estimator (QMPE) of $\boldsymbol{\theta}$, is given by $\hat{\boldsymbol{\theta}}_{\phi_{\lambda_2}} = \arg \min_{\boldsymbol{\theta} \in \Theta} d_{\phi_{\lambda_2}}(\hat{\mathbf{p}}, \mathbf{p}(\boldsymbol{\theta}))$, and the semi-parametric clustered overdispersed power-divergence based GOF test-statistic, based on $\hat{\boldsymbol{\theta}}_{\phi_{\lambda_2}}$, by

$$\tilde{T}_{\lambda_1, \lambda_2} = \frac{2\hat{n}N d_{\phi_{\lambda_1}}(\hat{\mathbf{p}}, \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_{\lambda_2}}))}{\tilde{\vartheta}_{\hat{n}^*, N, \phi_{\lambda_2}}} = \frac{2\hat{n}N}{\tilde{\vartheta}_{\hat{n}^*, N, \phi_{\lambda_2}} \lambda_1(\lambda_1 + 1)} \left(\sum_{r=1}^M \frac{\hat{p}_r^{\lambda_1+1}}{p_r^{\lambda_1}(\hat{\boldsymbol{\theta}}_{\phi_{\lambda_2}})} - 1 \right), \text{ for } \lambda_1 \notin \{-1, 0\}, \quad (4.1)$$

where $\tilde{\vartheta}_{\hat{n}^*, N, \phi_{\lambda_2}}$ is (3.3) and $\hat{n}N$ (3.5). The expression of the semiparametric clustered overdispersed power-divergence based GOF test-statistic for $\lambda_1 = 0$ ($\tilde{T}_{0, \lambda_2} = G^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_{\lambda_2}}) / \tilde{\vartheta}_{\hat{n}^*, N, \phi_{\lambda_2}}$) is in Corollary 3.4 and for the case of $\lambda_1 = -1$ is given by

$$\tilde{T}_{-1, \lambda_2} = \frac{\hat{n}N}{\tilde{\vartheta}_{\hat{n}^*, N, \phi_{\lambda_2}}} \sum_{r=1}^M p_r(\hat{\boldsymbol{\theta}}_{\phi_2}) \log \frac{p_r(\hat{\boldsymbol{\theta}}_{\phi_2})}{\hat{p}_r}.$$

Notice that the case of $\lambda_2 = 0$ for the QMPE of $\boldsymbol{\theta}$, matches the QMLE of $\boldsymbol{\theta}$, $\hat{\boldsymbol{\theta}}$, or equivalently the QM ϕ E of $\boldsymbol{\theta}$ with $\phi(x) = x \log x - x + 1$, and from the case of $\lambda_1 = 1$ arises the semiparametric clustered overdispersed chi-square GOF test-statistic $\tilde{T}_{1, \lambda_2} = X^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_{\lambda_2}}) / \tilde{\vartheta}_{\hat{n}^*, N, \phi_{\lambda_2}}$ given in Corollary 3.2. All these test-statistics are completely new when no distributional assumption is made, and homogeneous intraclass correlation assumption is considered cell by cell in all the clusters.

$\tilde{T}_{\lambda_1, \lambda_2}$ (p -value)		λ_2				
		-0.5	0	2/3	1	2
λ_1	-0.5	7.5621 (0.1090)	11.2413 (0.0240)	15.6963 (0.0035)	17.6234 (0.0015)	22.1483 (0.0002)
	0	7.7504 (0.1012)	9.7014 (0.0458)	12.2489 (0.0156)	13.4095 (0.0094)	16.2120 (0.0027)
	2/3	10.4138 (0.0340)	10.3330 (0.0352)	11.3428 (0.0230)	11.9922 (0.0174)	13.7789 (0.0080)
	1	13.0422 (0.0111)	11.2813 (0.0236)	11.4143 (0.0223)	11.8302 (0.0187)	13.2202 (0.0102)
	2	33.6045 (< 0.0001)	17.5637 (0.0015)	13.0587 (0.0110)	12.5518 (0.0137)	12.6781 (0.0130)
		2.1815	1.5869	1.3314	1.2707	1.1813
$\tilde{\vartheta}_{\hat{n}^*, N, \phi_{\lambda_2}}$		2.1815	1.5869	1.3314	1.2707	1.1813

Table 4.2: Values for the clustered overdispersed GOF test-statistic, via semi-parametric estimates of the design effect, with corresponding p -values.

The Brier's non-parametric estimator of ϑ_n can be also plugged on the clustered overdispersed GOF test-statistic,

$$\hat{T}_{\lambda_1, \lambda_2} = 2\hat{n}Nd_{\phi_{\lambda_1}}(\hat{\mathbf{p}}, \mathbf{p}(\hat{\boldsymbol{\theta}}_{\phi_{\lambda_2}}))/\hat{\vartheta}_{\hat{n}^*, N},$$

with no change in the asymptotic distribution. In particular,

$$\hat{T}_{0,0} = G^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_{\lambda_2}})/\hat{\vartheta}_{\hat{n}^*, N} \quad \text{and} \quad \hat{T}_{1,0} = X^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_{\lambda_2}})/\hat{\vartheta}_{\hat{n}^*, N}$$

are the clustered overdispersed GOF test-statistics proposed by Brier (1980).

$\hat{T}_{\lambda_1, \lambda_2}$ (p -value)		λ_2				
		-0.5	0	2/3	1	2
λ_1	-0.5	15.4857 (0.0038)	16.7462 (0.0022)	19.6173 (0.0006)	21.0219 (0.0003)	24.5600 (0.0001)
	0	15.8714 (0.0032)	14.4521 (0.0060)	15.3087 (0.0041)	15.9953 (0.0030)	17.9773 (0.0012)
	2/3	21.3256 (0.0003)	15.3931 (0.0040)	14.1762 (0.0068)	14.3048 (0.0064)	15.2792 (0.0042)
	1	26.7079 (< 0.0001)	16.8057 (0.0021)	14.2656 (0.0065)	14.1115 (0.0069)	14.6597 (0.0055)
	2	68.8157 (< 0.0001)	26.1646 (< 0.0001)	16.3207 (0.0026)	14.9723 (0.0048)	14.0586 (0.0071)
$\hat{\vartheta}_{\hat{n}^*, N}$		1.0653	1.0653	1.0653	1.0653	1.0653

Table 4.3: Values for the clustered overdispersed GOF test-statistic, via non-parametric estimates of the design effect, with corresponding p -values.

From the p -values of Tables 4.2 and 4.3 is concluded that only $\tilde{T}_{-0.5, -0.5}$ and $\tilde{T}_{0, -0.5}$ clustered overdispersed GOF test-statistics do not allow rejecting the null hypothesis.

5 Simulation Study

In the simulation study performed In Section 6.1 of Alonso-Revilla et al. (2016), a clear improvement of the semi-parametric estimator of ρ^2 , via QM ϕ Es, $\tilde{\rho}_{\hat{n}^*, N, \phi}^2$, was shown in comparison with the Brier's non-parametric estimator of ρ^2 , $\hat{\rho}_{\hat{n}^*, N}^2 = (\hat{\vartheta}_{\hat{n}^*, N} - 1)/(\hat{n}^* - 1)$. Taking into account the same simulation experiment, $\boldsymbol{\theta} = (\theta_{1(1)}, \theta_{1(2)}, \theta_{2(1)}, \theta_{2(2)})^T = (0.1, 0.2, 0.4, 0.3)^T$ is the true value of the parameter of the independence model described in Section 4, under the null hypothesis. The study considers $G = 3$ different cluster sizes with $N_1 = 18$, $N_2 = 2$, $N_3 = 5$ clusters, having each $n_1 = 5$, $n_2 = 3$, $n_3 = 7$ possibly correlated individuals.

With $R = 10,000$ replications the significance levels are estimated by simulation for the power divergence based GOF test-statistics $\tilde{T}_{\lambda_1, \lambda_2}$ and $\hat{T}_{\lambda_1, \lambda_2}$, with $\lambda_1, \lambda_2 \in \{-0.5, 0, 2/3, 1, 2\}$, defined in Section 4. An extensive study has been done by considering three possible distributions for $\mathbf{Y}^{(\ell)}$ but in Figure 1 only a summary of the final plots are shown. The three distributions, Dirichlet-multinomial (DM), random-clumped (RC) and n -inflated (NI), mentioned in Section 1, are generated according to the algorithms described in Alonso-Revilla et al. (2016) and Raim et al. (2015).

From the study it is concluded that a good behaviour of the estimator of ρ^2 (or ϑ_n) plays a crucial role on the behaviour of the closeness of the estimated significance level with respect to the nominal significance level, but the choice of $\lambda_1 = 2/3$ for the GOF test-statistic is also important. The combination of $\lambda_1 = 2/3$ for the GOF test-statistic with $\lambda_2 = 2$ for the estimator in $\tilde{T}_{\lambda_1, \lambda_2}$ (or $\lambda_1 = 2/3$ for the GOF test-statistic with $\lambda_2 = 0$ for the estimator) does not suffer negative modifications as the value of ρ^2 increases in the abscissa axis. The Brier's non-parametric estimator has however a negative impact on the estimated significance levels of the classical overdispersed likelihood-ratio GOF test-statistic $\hat{T}_{0,0} = G^2(\mathbf{Y}, \hat{\boldsymbol{\theta}}_{\phi_{\lambda_2}}) / \hat{\vartheta}_{\hat{n}^*, N}$ as the value of ρ^2 increases in the abscissa axis. Looking at the right hand side plots, the three distributions have estimated significance levels no closer to the nominal level, 0.05, in comparison with the rest of the distributions. In particular for $\lambda_1 = 2/3$ and $\lambda_2 = 2$ with the n -inflated distribution the estimated significance level tends to be below the nominal significance level, while for the Dirichlet-multinomial and random-clumped distribution, above the nominal significance level.

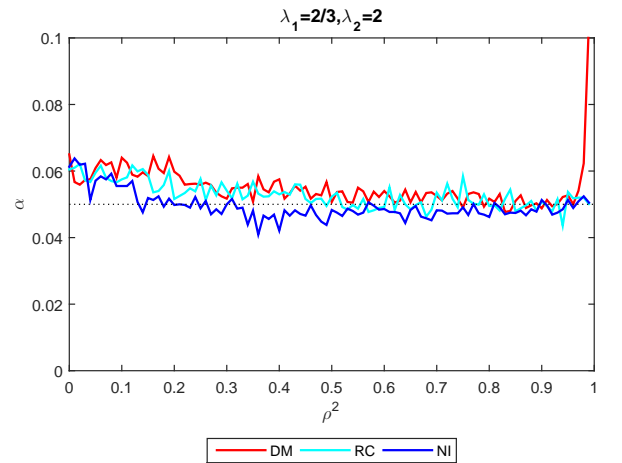
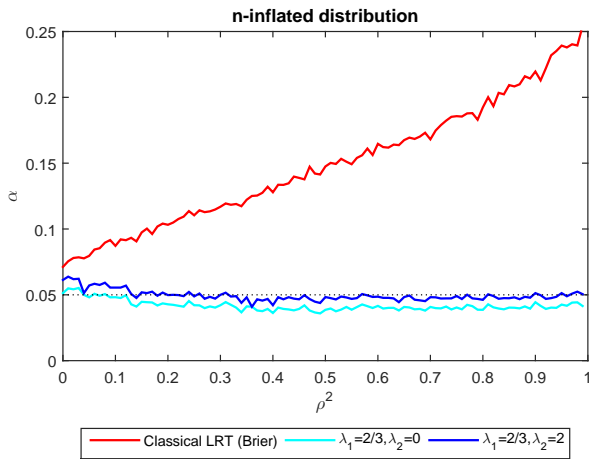
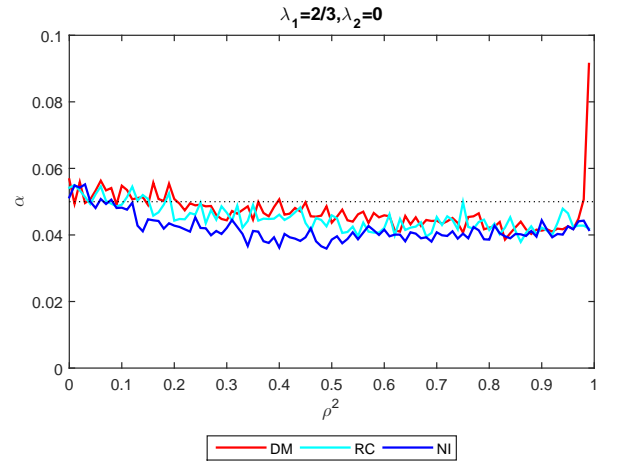
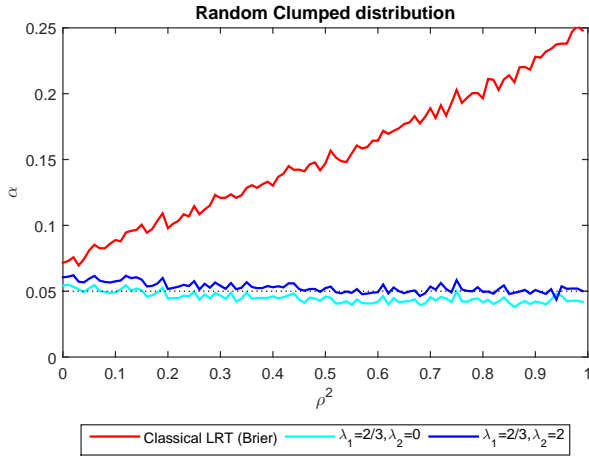
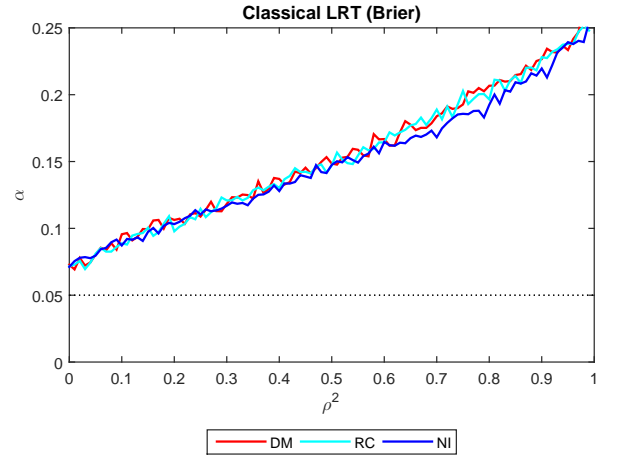
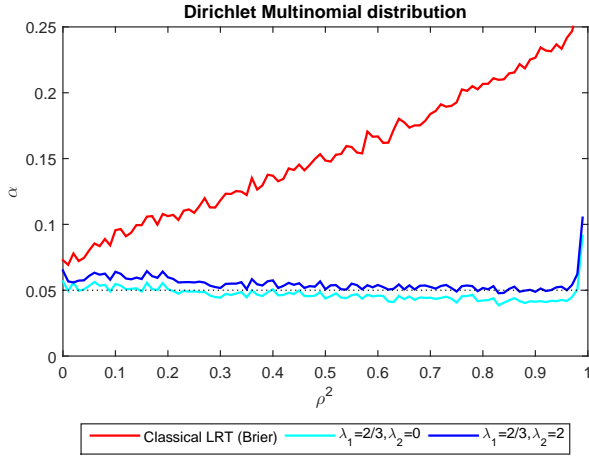


Figure 1: Estimated significance levels, by simulation, for three different distributions and types of overdispersed GOF test-statistics.

References

- [1] Alonso-Revenge, J. M., Martín, N. and Pardo, L. (2016). New improved estimators for overdispersion in models with clustered multinomial data and unequal cluster sizes. *Statistics and Computing* (in Press), DOI: 10.1007/s11222-015-9616-z.
- [2] Altham, P. M. E. (1976). Discrete variable analysis for individuals grouped into families. *Biometrika*, **63**, 263–269.
- [3] Bedrick, E. J. (1983). Adjusted chi-squared tests for cross-classified tables of survey data. *Biometrika*, **70**, 591–595.
- [4] Brier, S. S. (1980). Analysis of contingency tables under cluster sampling. *Biometrika*, **67**, 591–596.
- [5] Choi J.W. and McHugh R.B. (1989). A reduction factor in goodness-of-fit and independence tests for clustered and weighted observations. *Biometrics*, **45**, 979–996.
- [6] Cohen, J. E. (1976). The distribution of the chi-squared statistic under clustered sampling from contingency tables. *Journal of the American Statistical Association*, **71**, 665–670.
- [7] Cressie, N., Pardo, L. (2002). Phi-Divergence statistics. In: ElShaarawi, A. H., Piegorich, W. W., eds. *Encyclopedia of Environmetrics*. Vol. 3. New York, Wiley, pp. 1551–1555.
- [8] Eldridge, S. M., Ukoumunne, O. C. and Carlin, J. B. (2009). The Intra-Cluster Correlation Coefficient in Cluster Randomized Trials: A Review of Definitions. *International Statistical Review*, **77**, 378–394.
- [9] Fay, R. E. (1985). Complex samples. *Journal of the American Statistical Association*, **80**, 148–157.
- [10] Fellegi, I. P. (1980). Approximate tests of independence and goodness of fit based upon stratified multistage samples. *Journal of the American Statistical Association*, **75**, 261–268.
- [11] Fienberg, S. E. (1979). The use of chi-square statistics for categorical data problems. *Journal of the Royal Statistical Society, Series B*, **41**, 54–64.
- [12] Holt, D., Scott, A. J. and Ewings, P. O. (1980). Chi-squared tests with survey data. *Journal of the Royal Statistical Society, Series A*, **143**, 302–320.
- [13] Koch, G. G., Freeman, D. H., and Freeman, J. L. (1975). Strategies in the multivariate analysis of data from complex surveys. *International Statistical Review*, **43**, 59–78.
- [14] Landis, J. R., Lepkowski, J. M., Eklund, S. A. and Stehouwer, S. A. (1984). *A Statistical Methodology for Analyzing Data from a Complex Survey: The First National Health and Nutrition Examination Survey*. Series 2 , No. 92. Hyattsville, Maryland: The National Center for Health Statistics.
- [15] Menéndez, M. L., Morales, D., Pardo, L. and Vajda, I. (1995). Divergence-based estimation and testing of statistical models of classification. *Journal of Multivariate Analysis*, **54**, 329–354.

- [16] Menéndez, M. L., Morales, D., Pardo, L. and Vajda, I. (1996). About divergence-based goodness-of-fit tests in the dirichlet-multinomial model. *Communications in Statistics - Theory and Methods*, **25**, 1119–1133.
- [17] Morel, J.G. and Nagaraj, N.K. (1993). A finite mixture distribution for modelling multinomial extra variation. *Biometrika*, **80**, 363–371.
- [18] Mosimann, J. E. (1962). On the compound multinomial distributions, the multivariate β -distribution and correlation among proportions. *Biometrika*, **49**, 65–82.
- [19] Pardo, L. (2006). *Statistical inference based on divergence measures*. Chapman & Hall/CRC, Boca Raton.
- [20] Raim, A. M. , Neerchal, N. K. and Morel, J. G. (2015). Modeling overdispersion in *R*. Technical Report HPCI-2015-1 UMBCH High Performance Computing Facility, University of Maryland.
- [21] Rao, J. N. K. and Scott, A. J. (1981). The analysis of categorical data from complex sample surveys: Chi-squared tests for goodness of fit and independence in two-way tables. *Journal of the American Statistical Association*, **76**, 221–230.
- [22] Rao, J. N. K. and Scott, A. J. (1984). On chi-squared tests for multiway contingency tables with cell proportions estimated from survey data. *The Annals of Statistics*, **12**, 46–60.